



22883

PATENT TRADEMARK OFFICE

This application is submitted in the name of the following inventor:

1

2

3

InventorCitizenshipResidence (City and State)

4

Kleiman, Steven

United States

Los Altos, California

5

6

The assignee is Network Appliance, Inc., a corporation having an office at

7

2770 San Tomas Expressway, Santa Clara, CA 95051.

8

9

Title of the Invention

10

11

Highly Available File Servers

12

13

Background of the Invention

14

15

1. *Field of the Invention*

16

17

The invention relates to storage systems.

18

19

2. *Related Art*

20

21

Computer storage systems are used to record and retrieve data. In some

22

computer systems, storage systems communicate with a set of client devices, and provide

1 services for recording and retrieving data to those client devices. Because data storage is
2 important to many applications, it is desirable for the services and data provided by the
3 storage system to be available for service to the greatest degree possible. It is therefore
4 desirable to provide storage systems that can remain available for service even in the face
5 of component failures in the storage system.

6
7 One known technique for provide storage systems that can remain available
8 for service is to provide a plurality of redundant storage elements, with the property that
9 when a first storage element fails, a second storage element is available to provide the
10 services and the data otherwise provided by the first. Transfer of the function of provid-
11 ing services from the first to the second storage element is called "failover." The second
12 storage element maintains a copy of the data maintained by the first, so that failover can
13 proceed without substantial interruption.

14
15 A first known technique for achieving failover is to cause the second stor-
16 age element to copy all the operations of the first. Thus, each storage operation com-
17 pleted by the first storage element is also completed by the second. This first known
18 technique is subject to drawbacks: (1) It uses a substantial amount of processing power
19 at the second storage element duplicating efforts of the first, most of which is wasted. (2)
20 It slows the first storage element in confirming completion of operations, because the first
21 storage element waits for the second to also complete the same operations.

1 A second known technique for achieving failover is to identify a sequence
2 of checkpoints at which the first storage element is at a consistent and known state. On
3 failover, the second storage element can continue operation from the most recent check-
4 point. For example, the NFS (Network File System) protocol requires all write opera-
5 tions to be stored to disk before they are confirmed, so that confirmation of a write op-
6 eration indicates a stable file system configuration. This second known technique is
7 subject to drawbacks: (1) It slows the first storage element in performing write opera-
8 tions, because the first storage element waits for write operations to be completely stored
9 to disk. (2) It slows recovery on failover, because the second storage element addresses
10 any inconsistencies left by failure of the first between identified checkpoints.

11
12 Accordingly, it would be advantageous to provide a storage system, and a
13 method for operating a storage system, that efficiently uses all storage system elements,
14 quickly completes and confirms operations, and quickly recovers from failure of any
15 storage element. This advantage is achieved in an embodiment of the invention in which
16 the storage system implements frequent and rapid checkpoints, and in which the storage
17 system rapidly distributes duplicate commands for those operations between checkpoints
18 among its storage elements.

19

Summary of the Invention

The invention provides a storage system that is highly available even in the face of component failures in the storage system, and a method for operating that storage system. A first and a second file server each includes a file server request log for storing incoming file server requests. Both the first and second file servers have access to a common set of mass storage elements. Each incoming file server request is copied to both the first and second file servers; the first file server processes the file server request while the second file server maintains a copy in its file server request log. Each file server operates using a file system that maintains consistent state after each file server request. On failover, the second file server can perform those file server requests in its file server request log since the most recent consistent state.

In a second aspect of the invention, a file server system provides mirroring of one or more mass storage elements. Each incoming file server request is copied to both the first file server and the second file server. The first file server performs the file server requests to modify a primary set of mass storage elements, and also performs the same file server requests to modify a mirror set of mass storage elements. The mirror mass storage elements are disposed physically separately from the primary mass storage elements, such as at another site, and provide a resource in the event the entire primary set of mass storage elements is to be recovered.

Brief Description of the Drawings

Figure 1 shows a block diagram of a highly available file server system.

Figure 2 shows a block diagram of a file server in the file server system.

Figure 3 shows a process flow diagram of operation of the file server system.

Detailed Description of the Preferred Embodiment

In the following description, a preferred embodiment of the invention is described with regard to preferred process steps and data structures. However, those skilled in the art would recognize, after perusal of this application, that embodiments of the invention may be implemented using one or more general purpose processors (or special purpose processors adapted to the particular process steps and data structures) operating under program control, and that implementation of the preferred process steps and data structures described herein using such equipment would not require undue experimentation or further invention.

File Server Pair and Failover Operation

Figure 1 shows a block diagram of a highly available file server system.

A file server system 100 includes a pair of file servers 110, both coupled to a common set of mass storage devices 120. A first one of the file servers 110 is coupled to a first I/O bus 130 for controlling a first selected subset of the mass storage devices 120. Similarly, a second one of the file servers 110 is coupled to a second I/O bus 130 for controlling a second selected subset of the mass storage devices 120.

Although both file servers 110 are coupled to all of the common mass storage devices 120, only one file server 110 operates to control any one mass storage device 120 at any designated time. Thus, even though the mass storage devices 120 are each controllable by only one file server 110 at a time, each of the mass storage devices 120 remains available even if one of its two associated file servers 110 fails.

In a preferred embodiment, the file server system 100 includes a pair of such file servers 110; however, in alternative embodiments, more than two such file servers 110 may be included in a single file server system 100.

In a preferred embodiment, the first I/O bus 130 and the second I/O bus 130 each include a mezzanine bus such as the PCI bus architecture.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

In a preferred embodiment, the mass storage devices 120 include magnetic disk drives, optical disk drives, or magneto-optical disk drives. In alternative embodiments, however, other storage systems may be used, such as bubble memory, flash memory, or systems using other storage technologies. Components of the mass storage devices 120 are referred to as "disks," even though those components may comprise other forms or shapes.

Each mass storage device 120 can include a single disk or a plurality of disks. In a preferred embodiment, each mass storage device 120 includes a plurality of disks and is disposed and operated as a RAID storage system.

In a preferred embodiment, the first file server 110 is coupled to the second file server 110 using a common interconnect. The common interconnect provides a remote memory access capability for each file server 110, so that data can be stored at each file server 110 from a remote location. In a preferred embodiment, the common interconnect includes a Tandem "ServerNet" interconnect. The common interconnect is coupled to each file server 110 using a device controller coupled to an I/O bus for each file server 110.

The first file server 110 is coupled to a first network interface 140, which is disposed to receive file server requests 151 from a network 150. Similarly, the second

1 file server 110 is coupled to a second network interface 140, which is also disposed to
2 receive file server requests 151 from the network 150.

3

4 The first file server 110 includes a first server request memory 160, which
5 receives the file server requests 151 and records them. In the event the first file server
6 110 recovers from a power failure or other service disruption, the outstanding file server
7 requests 151 in the first server request memory 160 are re-performed to incorporate them
8 into a next consistent state of the file system maintained by the first file server 110.

9

10 Similarly, the second file server 110 includes a second server request mem-
11 ory 160, which receives the file server requests 151 and records them. In the event the
12 second file server 110 recovers from a power failure or other service disruption, the out-
13 standing file server requests 151 in the second server request memory 160 are re-
14 performed to incorporate them into a next consistent state of the file system maintained
15 by the second file server 110.

16

17 When the first file server 110 receives a file server request 151 from the
18 network 150, that file server request 151 is copied into the first server request memory
19 160. The file server request 151 is also copied into the second server request memory
20 160 using remote memory access over the common interconnect. Similarly, when the
21 second file server 110 receives a file server request 151 from the network 150, that file
22 server request 151 is copied into the second server request memory 160. The file server

1 request 151 is also copied into the first server request memory 160 using remote memory
2 access over the common interconnect. Using remote memory access is relatively quicker
3 and has less communication overhead than using a networking protocol.

4

5 In the event that either file server 110 fails, the other file server 110 can
6 continue processing using the file server requests 151 stored in its own server request
7 memory 160.

8

9 In a preferred embodiment, each server request memory 160 includes a
10 nonvolatile memory, so those file server requests stored in either server request memory
11 160 are not lost due to power failures or other service interruptions.

12

13 The responding file server 110 processes the file server request 151 and
14 possibly modifies stored files on one of the mass storage devices 120. The non-
15 responding file server 110, partner to the responding file server 110, maintains the file
16 server request 151 stored in its server request memory 160 to prepare for the possibility
17 that the responding file server 110 might fail. In the event the responding file server 110
18 fails, the non-responding file server 110 processes the file server request 151 as part of a
19 failover technique.

20

1 In a preferred embodiment, each file server 110 controls its associated mass
2 storage devices 120 so as to form a redundant array, such as a RAID storage system, us-
3 ing inventions described in the following patent applications:

4
5 o Application Serial No. 08/471,218, filed June 5, 1995, in the name of inventors
6 David Hitz et al., titled "A Method for Providing Parity in a Raid Sub-System
7 Using Non-Volatile Memory", attorney docket number NET-004;

8
9 o Application Serial No. 08/454,921, filed May 31, 1995, in the name of inventors
10 David Hitz et al., titled "Write Anywhere File-System Layout", attorney docket
11 number NET-005;

12
13 o Application Serial No. 08/464,591, filed May 31, 1995, in the name of inventors
14 David Hitz et al., titled "Method for Allocating Files in a File System Integrated
15 with a Raid Disk Sub-System", attorney docket number NET-006.

16
17 Each of these applications is hereby incorporated by reference as if fully set
18 forth herein. They are collectively referred to as the "WAFL Disclosures."

19
20 As part of the techniques shown in the WAFL Disclosures, each file server
21 110 controls its associated mass storage devices 120 in response to file server requests
22 151 in an atomic manner. The final action for any file server request 151 is to incorpo-

1 rate the most recent consistent state into the file system 121. Thus, file system 121 is in
2 an internally consistent state after completion of each file server request 151. Thus, a file
3 system 121 defined over the mass storage devices 120 will be found in an internally con-
4 sistent state, regardless of which file server 110 controls those mass storage devices 120.
5 Exceptions to the internally consistent state will only include a few of the most recent file
6 server requests 151, which will still be stored in the server request memory 160 for both
7 file servers 110. Those most recent file server requests 151 can be incorporated into a
8 consistent state by performing them with regard to the most recent consistent state.

9
10 For any file server request 151, in the event the file server 110 normally re-
11 sponding to that file server request 151 fails, the other file server 110 will recognize the
12 failure and perform a failover method to take control of mass storage devices 120 previ-
13 ously assigned to the failing file server 110. The failover file server 110 will find those
14 mass storage devices 120 with their file system 121 in an internally consistent state, but
15 with the few most recent file server requests 151 as yet unperformed. The failover file
16 server 110 will have copies of these most recent file server requests 151 in its server re-
17 quest memory 160, and will perform these file server requests 151 in response to those
18 copies.

19

20 *File Server Node*

21

22

Figure 2 shows a block diagram of a file server in the file server system.

Each file server 110 includes at least one processor 111, a program and data memory 112, the server request memory 160 (including a nonvolatile RAM), a network interface element 114, and a disk interface element 115. These elements are interconnected using a bus 117 or other known system architecture for communication among processors, memory, and peripherals.

In a preferred embodiment, the network interface element 114 includes a known network interface for operating with the network 150. For example, the network interface element 114 can include an interface for operating with the FDDI interface standard or the 100BaseT interface standard.

After failover, the file server 110 responds to file server requests directed to either itself or its (failed) partner file server 110. Each file server 110 is therefore capable of assuming an additional network identity on failover, one for itself and one for its failed partner file server 110. In a preferred embodiment, the network interface element 114 for each file server 110 includes a network adapter capable of responding to two separate addresses upon instruction by the file server 110. In an alternative embodiment, each file server 110 may have two such network adapters.

In a preferred embodiment, the disk interface element 115 includes a known disk interface for operating with magnetic, optical, or magneto-optical disks, that

has two independent ports with each port coupled to a separate file server 110, such as the FC-AL interface. This helps prevent failure of one file server 110 from affecting low-level operation of the other file server 110.

In a preferred embodiment, the bus 117 includes at least a memory bus 171 and the mezzanine bus 130. The memory bus 171 couples the processor 111 and the program and data memory 112. The mezzanine bus 130 couples the network interface element 114 and the disk interface element 115. The memory bus 171 is coupled to the mezzanine bus 130 using an I/O controller 173 or other known bus adapter technique.

In a preferred embodiment, each disk in the mass storage 120 is statically assigned to either the first file server 110 or the second file server 110, responsive to whether the disk is wired for primary control by either the first file server 110 or the second file server 110. Each disk has two control ports A and B; the file server 110 wired to port A has primary control of that disk, while the other file server 110 only has control of that disk when the other file server 110 has failed.

Operation Process Flow

Figure 3 shows a process flow diagram of operation of the file server system.

1 A method 300 is performed by the components of the file server 100, and
2 includes a set of flow points and process steps as described herein.

3
4 At a flow point 310, a device coupled to the network 150 desires to make a
5 file system request 151.

6
7 At a step 311, the device transmits the file system request 151 to the net-
8 work 150.

9
10 At a step 312, the network 150 transmits the file server request 151 to the
11 file server 110.

12
13 At a step 313, a first file server 110 at the file server system 100 receives
14 the file server request 151. The first file server 110 copies the file server request 151 into
15 the first server request memory 160, and also copies the file server request 151 into the
16 second server request memory 160 using the common interconnect. The target of the
17 copying operation in the second server request memory 160 is to an area reserved for this
18 purpose. The copying operation requires no further processing by the second file server
19 110, and the second file server 110 does not normally process or respond to the file server
20 request 151.

1 At a step 314, the first file server 110 responds to the file server request
2 151.

3
4 At a flow point 320, the file server request has been successfully processed.
5

6 In a second aspect of the invention, the first file server 110 provides mir-
7 roring of one or more of its mass storage devices 120.
8

9 As with the first aspect of the invention, each incoming file server request
10 is copied to both the first file server 110 and the second file server 110. The first file
11 server 110 performs the file server requests to modify one or more primary mass storage
12 devices 120 under its control. The first file server 110 also performs the file server re-
13 quests to modify a set of mirror mass storage devices 120 under its control, but located
14 distant from the primary mass storage devices 120. Thus, the mirror mass storage de-
15 vices 120 will be a substantial copy of the primary mass storage devices 120.
16

17 The mirror set of mass storage devices 120 provide a resource in the event
18 the entire primary set of mass storage devices 120 is to be recovered, such as if a disaster
19 befalls the primary set of mass storage devices 120.
20

21 At a flow point 330, the first file server 110 in the file server system 100
22 fails.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

At a step 331, the second file server 110 in the file server system 100 recognizes the failure of the first file server 110.

In a preferred embodiment, the second file server 110 performs the step 331 in the following manner:

o Each file server 110 maintains two disks of its mass storage devices 120 (thus, there are a total of four such disks for two file servers 110) for recording state information about the file server 110. There are two such disks (called "mailbox disks") so that one can be used as primary storage and one can be used as backup storage. If one of the two mailbox disks fails, the file server 110 using that mailbox disk designates another disk as one of its two mailbox disks.

o Each file server 110 maintains at least one sector on each mailbox disk, on which the file server 110 periodically writes state information. Each file server 110 also sends its state information to the other file server 110 using the interconnect using remote memory access. The state information written to the mailbox disks by each file server 110 changes with each update.

o Each file server 110 periodically reads the state information from at least one of the mailbox disks for the other file server 110. Each file server 110 also receives

1 state information from the other file server 110 using the interconnect using re-
2 mote memory access.

3
4 o Each file server 110 recognizes if the other file server 110 has failed by noting that
5 there has been no update to the state information on the mailbox disks for the other
6 file server 110.

7
8 In a preferred embodiment, the second file server 110 determines whether
9 failure of the first file server 110 is a hardware error or a software error, and only recog-
10 nizes failure of the first file server 110 for hardware errors. In alternative embodiments,
11 the second file server 110 may recognize failure of the first file server 110 for software
12 errors as well.

13
14 At a step 332, the second file server 110 seizes control of all mass storage
15 devices 120 previously assigned to the first file server 110. Due to the nature of the tech-
16 niques shown in the WAFL Disclosures, the file system 121 defined over those mass
17 storage devices 120 will be in an internally consistent state. All those file server requests
18 151 marked completed will have been processed and the results incorporated into storage
19 blocks of the mass storage devices 120.

20
21 In normal operation, neither file server 110 places reservations on any of
22 the mass storage devices 120. In the step 332 (only on failover), the second file server

1 110 seizes control of the mass storage devices 120 previously controlled by the first file
2 server 110, and retains control of those mass storage devices 120 until it is satisfied that
3 the first file server 110 has recovered.

4

5 When the first file server 110 recovers, it sends a recovery message to the
6 second file server 110. In a preferred embodiment, the second file server 110 relin-
7 quishes control of the seized mass storage devices 120 by operator command. However,
8 in alternative embodiments, the second file server 110 may recognize the recovery mes-
9 sage from the first file server 110 and relinquish control of the seized mass storage de-
10 vices 120 in response thereto.

11

12 At a step 333, the second file server 110 notes all file server requests 151 in
13 the area of its server request memory 160 that were copied there by the first file server
14 110. Those file server requests 151 whose results were already incorporated into storage
15 blocks of the storage devices 120 are discarded.

16

17 At a step 334, when the second file server 110 reaches its copy of each file
18 server request 151, the second file server 110 processes the file server request 151 nor-
19 mally.

20

21 At a flow point 340, failover from the first file server 110 to the second file
22 server 110 has been successfully handled.

1

2 *Alternative Embodiments*

3

4 Although preferred embodiments are disclosed herein, many variations are
5 possible which remain within the concept, scope, and spirit of the invention, and these
6 variations would become clear to those skilled in the art after perusal of this application.